

AD-A152 136

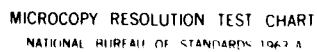
IMPLEMENTATION OF KOREAN AND CHINESE CHARACTERS THROUGH 1/1
COMPUTER(U) NAVAL POSTGRADUATE SCHOOL MONTEREY CA
C H KIM ET AL. SEP 84

UNCLASSIFIED

F/G 9/2

NL

										END			
										FILED			
										DATE			



2

NAVAL POSTGRADUATE SCHOOL

Monterey, California

AD-A152 136



DTIC
ELECTE
S APR 5 1985 D
A

THESIS

IMPLEMENTATION OF KOREAN AND CHINESE
CHARACTERS THROUGH COMPUTER

by

Chong Hae Kim
and
Sung Woo Ko

September 1984

Thesis Advisor:

Michael J. Zyda

Approved for public release; distribution unlimited

DTIC FILE COPY

85 03 18 006

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER	2. GOVT ACCESSION NO. -10-A/52136	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) Implementation of Korean and Chinese Characters through Computer		5. TYPE OF REPORT & PERIOD COVERED Master's Thesis September 1984
		6. PERFORMING ORG. REPORT NUMBER
7. AUTHOR(s) Chong Hae Kim and Sung Woo Ko		8. CONTRACT OR GRANT NUMBER(s)
9. PERFORMING ORGANIZATION NAME AND ADDRESS Naval Postgraduate School Monterey, California 93943		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
11. CONTROLLING OFFICE NAME AND ADDRESS Naval Postgraduate School Monterey, California 93943		12. REPORT DATE September 1984
		13. NUMBER OF PAGES 62
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		15. SECURITY CLASS. (of this report) UNCLASSIFIED
		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) bit, byte, code, keyboard, keytop, bit-map, matrix printer, CRT, resolution, quality, enumeration.		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) Methods of representing Korean and Chinese characters are presented, using a limited number of keystrokes on a standard keyboard. Various attempts have been made to find the most efficient way to represent these characters such as enumeration methods, 16-bit coding for Korean character syllables, and the meaning and the sound method for Chinese characters. Details of these are explained with a brief introduction (Continued)		

ABSTRACT (Continued)

to some general properties of Korean and Chinese characters currently used in Korea.

Accession For	
NTIS GRA&I	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	<input type="checkbox"/>
By _____	
Distribution/	
Availability	
Distribution	
Dist _____	
H-1	



Approved for public release; distribution unlimited.

Implementation of Korean and Chinese Characters
through Computer

by

Chong Hae Kim
Major, Republic of Korea Army
B.S., Korea Military Academy, 1976

and

Sung Woo Ko
Major, Republic of Korea Army
B.A., Korea Military Academy, 1976

Submitted in partial fulfillment of the
requirements for the degree of


MASTER OF SCIENCE IN INFORMATION SYSTEMS

from the

NAVAL POSTGRADUATE SCHOOL
September 1984

Authors:

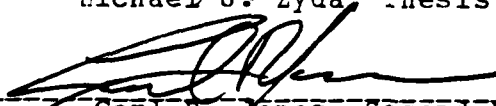

Chong Hae Kim


Sung Woo Ko

Approved by:



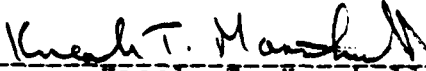
Michael J. Zyda, Thesis Advisor



Carl R. Jones, Second Reader



Willis H. Greer, Jr., Chairman,
Department of Administrative Sciences



Kneale T. Marshall,
Dean of Information and Policy Sciences

ABSTRACT

Methods of representing Korean and Chinese characters are presented, using a limited number of keystrokes on a standard keyboard. Various attempts have been made to find the most efficient way to represent these characters such as enumeration methods, 16-bit coding for Korean character syllables, and the meaning and the sound method for Chinese characters. Details of these are explained with a brief introduction to some general properties of Korean and Chinese characters currently used in Korea.

TABLE OF CONTENTS

I.	INTRODUCTION	9
II.	BACKGROUND	10
	A. PROPERTIES OF KOREAN DOCUMENTS	10
	B. CHARACTERISTICS OF KOREAN SCRIPT	10
	C. CHARACTERISTICS OF SINO-KOREAN CHARACTERS	14
III.	PROBLEMS OF EDITING KOREAN AND CHINESE SCRIPTS	17
	A. CURRENT EDITING TECHNOLOGY	17
	B. USER REQUIREMENTS	18
	C. PROBLEMS OF REPRESENTATION OF THE THREE KINDS OF SCRIPTS	19
IV.	POSSIBLE METHODS FOR KOREAN LANGUAGE DATA PROCESSING	20
	A. 8-BIT CODE FOR KOREAN ALPHABET	20
	1. Using the Current Standard Keyboard	20
	2. Using the Capital Letters as the Initial Letter	21
	B. 16-BIT CODE FOR THE THREE KINDS OF SCRIPT	22
	1. 16-bit Code for Korean Script	22
	2. 16-bit Code for Chinese Characters	25
	3. 16-bit Code for Roman Alphabet and Symbols	32
	4. Keyboard for 16-bit Code	32
	5. Operating System for Input and Output	35
	6. Design Considerations for Character Generation	37
	7. Memory Space for Character Definition	42

V.	EVALUATION OF SUGGESTED METHODS	45
VI.	RECCMMENDATION AND CONCIUSION	49
APPENDIX A:	THE EVOLUTION OF CHINESE CHARACTERS	51
APPENDIX B:	EBCDIC INPUT CODE	52
APPENDIX C:	MDS INPUT CODE	53
APPENDIX D:	IBM MULTISTATION 5550 KEYBOARD	54
APPENDIX E:	FACOM OS IV (KEF) KEYBOARD	55
APPENDIX F:	LOAD COMMAND PROGRAM FOR CURRENT KEYBOARD	56
APPENDIX G:	LOAD COMMAND PROGRAM FOR CAPITAL LETTER KEYBOARD	57
APPENDIX H:	IBM 2-BYTE INTERNAL HANGUL CODE	58
LIST OF REFERENCES	59
BIBLICGRAPHY	61
INITIAL DISTRIBUTION LIST	62

LIST OF TABLES

I.	Proportions of Written Characters	11
II.	Frequency of Chinese Characters Used in Documents	15
III.	Structure of 16-bit Code for Korean Script	23
IV.	16-bit Code for Korean Script	24
V.	Structure of 16-bit Code for Chinese Characters	26
VI.	Characters Having Same 1st & 2nd Letter Sound	27
VII.	5-bit Code for the Acronym of Sound Character	28
VIII.	Frequency of Meaning Character in Chinese Characters	30
IX.	5-bit Code for Acronym of Meaning Character	31
X.	Proportion of Duplicate Code	32
XI.	Structure of 16-bit Code for Roman Alphabet	33
XII.	Leveled Letters on 32 Key Tops	34
XIII.	Typing Procedures for Mixed Characters	35
XIV.	Comparative Table for Performance	40
XV.	Varieties of Data Compression Methods	43
XVI.	Net Present Value Formula	47

LIST OF FIGURES

2.1	The Korean Alphabet	13
4.1	Example Using Standard Keyboard	21
4.2	Example Using Capital Letters	22
4.3	Flowchart of Input and Output Controller	36
4.4	An Example of Gothic Type	38
4.5	An Example of Brush Type	38
4.6	Criteria in Designing a Character Generator	39

I. INTRODUCTION

The development of computer and information processing has come to the stage of being able to handle Korean and Chinese character input and output. There is no problem in information systems for the input and output of characters from a standard Roman character keyboard, but the problems related to non-Roman characters from I/O to software problems of language handling remain almost unsolved. Until recently the computer could not handle Korean or Chinese characters efficiently. It was not user friendly and data processing in Korea was imperfect and very unwieldy. Among the problems, the biggest issue is how to enter 2,369 Korean and 1,800 common Chinese characters from the standard Roman character keyboard.

During the last few years, there have been great efforts at universities, research institutes and manufacturers for the development of good I/O devices for Korean characters. In Korea, natural language processing, especially Korean language processing, is one of the essential elements for the future of computer and information systems.

First the properties of Korean and Chinese characters will be presented as an introduction for those unfamiliar with these characters. Then, the resolution power of CRT's and dot matrix printers and their relation to the shape characteristics (readability, aesthetic quality, etc.) of Korean and Chinese characters will be discussed. The methods which are developed for Korean and Chinese character I/O can be applied to other character sets, especially to many non-Roman alphabetic character sets, not to mention Chinese characters in China.

II. BACKGROUND

A. PROPERTIES OF KOREAN DOCUMENTS

Common documents in Korea are usually written in a mixed form utilizing Korean and Chinese characters. Minor use is made of Roman script. The usage of each character set depends on the kind of document. In order to perform word processing efficiently in Korea, the simultaneous editing of these characters is essential. Table I shows the use of characters found for various types of documents. This data is based on sampling performed expressly for this study. The following sources were in the sampling process to construct Table I:

1. newspaper - Korean Daily Times, "3A Era", 16 September 1984
2. journal - "National Security", June 1984
3. technical papers (A) - "COBOL Programming", Dong-A publishing Co., 1978
4. technical papers (B) - "Introduction to Law", Beoh Moon Sa publishing Co., 1978
5. business papers - Korean Air Lines Co.

Although the sample was taken from a single source for each kind of document, it is the authors' view that the documents selected are representative of the entire population of each type.

B. CHARACTERISTICS OF KOREAN SCRIPT

The native Korean alphabet was introduced in 1446, after centuries of the use of a more cumbersome method (known as IDU) to transcribe Korean with Chinese characters. The set

TABLE I
Proportions of Written Characters

	News- paper	Journal	Technical paper (A)	paper (B)	Business paper
Roman script	1%	3%	40%	0%	10%
Korean script	84%	76%	55%	55%	80%
Chinese character	15%	21%	5%	45%	10%

* (A): Technical papers from western countries
* (B): Traditional and historical papers

of 28 letters¹ (now 24 letters) was designed by a group of scholars commissioned by King Sejong (1419 - 1450), the fourth King of the Yi dynasty.

The Korean language and alphabet is spoken and written by an estimated 50 million people on the Korean peninsula and its coastal islands. Many among the approximately one million Koreans residing in Japan, China, and America still speak and write the language [Ref. 9].

The Korean alphabet currently used consists of 14 consonants (ㄱ ㄴ ㄷ ㄹ ㅁ ㅂ ㅅ ㅇ ㅈ ㅊ ㅋ ㅌ ㅍ ㅎ) and 10 vowels (ㅏ ㅑ ㅓ ㅕ ㅗ ㅛ ㅜ ㅠ ㅡ ㅣ). There are also 17 compound consonants (ㄱㅈ ㄴㅈ ㄴㅊ ㄷㅌ ㄹㅌ ㄹㅍ ㄹㅎ ㅁㅂ ㅁㅅ ㅅㅈ ㅈㅈ) and 11 compound vowels (ㅏㅑ ㅑㅓ ㅓㅕ ㅕㅗ ㅗㅛ ㅛㅜ ㅜㅠ ㅠㅡ ㅡㅣ ㅏㅓ ㅓㅕ). The letters of the Korean alphabet cannot be used independently but are used to build syllables. Each Korean character consists of two or three parts. The first part

¹A letter is an element of a character. The character consists of two or three letters. Letters in Korea are a set of 14 consonants and 10 vowels.

- * . THE CHARACTERISTICS OF KOREAN CHARACTER
 - . THE KOREAN ALPHABET CONSISTS OF 24 BASIC LETTERS(ELEMENTS);
 - 14 CONSONANTS: ㄱ ㄴ ㄷ ㄹ ㅁ ㅂ ㅅ ㅈ ㅊ ㅋ ㆁ ㄷ ㄷ ㄷ
 - 10 VOWELS : ㅏ ㅑ ㅓ ㅕ ㅗ ㅛ ㅜ ㅠ ㅡ ㅣ
 - . EACH CONSONANT AND VOWEL CAN BE COMPOUNDED
 - . POSSIBLE COMPOUND CONSONANTS
 - ㄱ ㄱ ㄴ ㄴ ㄷ ㄷ ㄹ ㄹ ㅁ ㅁ ㅂ ㅂ ㅅ ㅅ ㅈ ㅈ ㅊ ㅊ ㅋ ㅋ ㆁ ㆁ
 - . POSSIBLE COMPOUND VOWELS
 - ㅏ ㅑ ㅓ ㅕ ㅗ ㅛ ㅜ ㅠ ㅡ ㅣ
- *. EACH CHARACTER CAN BE DIVIDED INTO THREE PARTS (FIRST SOUND, MIDDLE SOUND, FINAL SOUND) OR TWO PARTS (FIRST AND SECOND SOUND).
 - . THE FIRST PART MUST CONSIST OF A CONSONANT OR A COMPOUND CONSONANT
THE SECOND PART MUST CONSIST OF A SINGLE OR A COMPOUND VOWEL
THE THIRD PART IS OPTIONAL. IF USED, IT MUST BE A CONSONANT.
 - . THE FOLLOWING LETTERS CAN BE USED AS THE FIRST PART;
 - ㄱ ㄴ ㄷ ㄹ ㅁ ㅂ ㅅ ㅈ ㅊ ㅋ ㆁ ㄷ ㄷ ㄷ
 - 19 LETTERS
 - . THE FOLLOWING LETTERS CAN BE USED AS THE SECOND PART;
 - ㅏ ㅑ ㅓ ㅕ ㅗ ㅛ ㅜ ㅠ ㅡ ㅣ
 - 21 LETTERS
 - . THE FOLLOWING LETTERS CAN BE USED AS THE THIRD PART;
 - ㄱ ㄴ ㄷ ㄹ ㅁ ㅂ ㅅ ㅈ ㅊ ㅋ ㆁ ㄷ ㄷ ㄷ
 - 28 LETTERS
- *. NUMBER OF POSSIBLE COMBINATIONS OF CHARACTER = $19 \times 21 \times 29 = 11,571$
IN PRACTICE, ONLY ABOUT 2,400 CHARACTERS ARE USED.

Figure 2.1 The Korean Alphabet.

top-to-bottom arrangement (e.g., ㅏ). The particular vowel being written determines which arrangement is used.

Representing these character syllables through a computer creates a problem because each letter's (consonant and vowel) shape can be different due to a requirement that each character be balanced, i.e., have the same size and achieve a desired aesthetic quality. For example, when ㄱ is placed to the left of a vowel, the downward portion is slanted: ㄱㅏ (e.g., ㅏ). When it is placed on top of the vowel, the downward portion becomes straight: ㄱㅏ (e.g., ㅏ). As shown above, it is very difficult to apply these different shapes for a particular letter to a line printer and a typewriter. This problem will be discussed in detail in the following chapter.

By mathematical calculation, the possible number of Korean characters is 11,571 (19 * 21 * 29). It must be noted through that only 2,369 characters are commonly used [Ref. 8: p. 11].

C. CHARACTERISTICS OF SINO-KOREAN CHARACTERS

Sino-Korean characters are Chinese characters used in Korea. They are different from those used in China. Koreans refer to Chinese characters as Hanja. Chinese characters have a long history, the earliest discovered writings having been dated from about 14 B.C.. In 100 A.D. during the Han Dynasty, this was modified by Hsu Sheng (許慎, 30 - 124) in his 15 - Volume paleographical work, Shuo-wen Chieh-tzu, (說文解字) which translates to the explanation of writing and analysis of words. That work lists 9,353 characters under 540 radical entries. Of this number, 364 are pictographic, 125 simple ideographic, 1,167 compound ideographic and 7,697 phonetic compounds.

The most complete collection, the Kang Hsi Dictionary with about 50,000 characters was published in 1716. Since 1949, after the establishment of the Peoples Republic of China, the Chinese government actively pursued language reform until the Cultural Revolution, 1966-1976. The Chinese government changed and simplified the characters from the original [Ref. 5: p. 15].

The number of characters used commonly is from 1,000 to 3,000. Table II [Ref. 1: p. 819] shows the frequency of

TABLE II
Frequency of Chinese Characters Used in Documents

	News- papers (%)	General Document (%)	Total Document (%)	News- papers (Chrs)	General Document (Chrs)
1st 10 chrs	10.0	8.8	80	499	638
50	27.5	25.5	85	615	777
100	38.9	36.1	90	781	992
200	55.4	51.0	95	1068	1358
500	79.0	73.5	96	1156	1479
1000	93.1	89.0	97	1269	1617
1500	97.4	95.0	98	1421	1832
2000	98.7	97.6	99	1661	2157
2500	98.9	99.4	100	2879	3328
3000		99.8			

* chrs: acronym of characters

Chinese characters used in typical documents.

In 1972, the Korean ministry of Education suggested that 1,800 Chinese characters be learned and used for educational purposes [Ref. 3]. In this study, the authors will restrict themselves to that set of 1,800 characters. The Chinese characters are called Hantzu in Chinese, Hanja in Korean, and Kanji in Japanese. All mean "Han Characters" (漢字).

These characters are used exclusively in Chinese writings, and in combination with the Hangul (Korean) alphabet in Korea and with the Kana Syllabaries in Japan. The Sino-Korean (Hanja), in written form, is a combination of three major elements: pictograms and ideograms, and phonograms [Ref. 5: p. 22].

In the next chapter the perspective of a picture for each character will be used because of both the complexity of Chinese characters and the ease of representation in the computer. Each Chinese character has the meaning and sound, for example, 天 means heaven and the sound is cheon. Also, there are many characters which have different meanings but the same sound, or the same meaning but different sounds. In order to solve this problem there are several methods. Appendix A [Ref. 5: p. 17] shows the evolution of Chinese characters.

III. PROBLEMS OF EDITING KOREAN AND CHINESE SCRIPTS

A. CURRENT EDITING TECHNOLOGY

The current word processing practice in Korea is to type Korean characters by the enumeration method, that is, input letters (8 bit code: consonant and vowel in sequence) and output these letters as a character syllable using a Korean character conversion program for Korean script. Appendix B shows the EBCDIC input codes currently used by FACOM, and Appendix C depicts MDS (Mahawk Data Sciences) input codes used by IBM. To type Chinese characters the following sequence is followed:

1. Depressing a Chinese character function key.
2. Typing the sound character of a Chinese character using the enumeration method.
3. Displaying all homonym (from 1 to 60) characters [Ref. 4: p. 34] that have the same sound.
4. Selecting one character by using an index number, and entering the character to a buffer or file.

Machines dealing with Korean language data are currently available from the IBM and FACOM corporations in Korea; IBM's Multistation 5550 (1984) and FACOM OS IV(KEF) (1982) are newly updated and well developed machines. These machines still have several disadvantages in handling Korean and Chinese characters:

1. A large amount of time is spent in character conversion.
2. It is difficult to directly delete and insert records in a file.
3. The word processing editor cannot recognize the characters being edited before executing a character

conversion program since only the enumerated letters can be displayed.

4. The method of entering characters is inconvenient and requires a tremendous amount of effort for Chinese characters.
5. One cannot convert all Korean character syllables into Chinese characters because there is not a one to one mapping.
6. Data communication is impossible since there are no standard codes for Korean and Chinese characters.

Appendices D and E show the keyboard of IBM Multistation 5550 [Ref. 8: p. 14] and FACOM OS IV (KEF) [Ref. 7: p. 48] respectively.

B. USER REQUIREMENTS

Most potential users have recognized that the computer is essential in data processing and office automation. However, because of the above constraints, they are unsatisfactory for use with the Korean language. Some general user requirements of computer researchers and manufacturers are the following:

1. Users want to use Korean language commands and programs but there are no Korean language oriented operating systems or programming languages such as COBOL, FORTRAN, Pascal, etc.
2. Users want to edit three kinds of characters simultaneously and in a user friendly manner.
3. Users want to display and print out data without using a conversion program, as is done with the Korean alphabet because of time, memory space, and inconvenience.
4. Users want to use interactive files and database processing.

In summation, they want to use computers that handle three kinds of script in the same manner in which present computers do with the Roman alphabet.

C. PROBLEMS OF REPRESENTATION OF THE THREE KINDS OF SCRIPTS

Because of the characteristics of Korean and Chinese characters, the following problems occur:

1. How can one enter 2,400 Korean characters and 1,800 Chinese characters into a computer through a limited number of keystrokes.
2. How can one develop the system program to direct input and output without using a conversion program.
3. How can the asthetic quality of display and output be improved.
4. How can one increase the processing speed and reduce the memory space for these character definitions.

There are other problems but the above problems are the most significant. Among these problems the first one is the most serious and significant problem, and consequently, the authors will give it more attention in this study.

IV. POSSIBLE METHODS FOR KOREAN LANGUAGE DATA PROCESSING

In order to solve the problems which were mentioned in the previous chapter, the following methods are offered as possible alternatives for Korean language data processing.

A. 8-BIT CODE FOR KOREAN ALPHABET

Since the Korean alphabet consists of only 24 letters and Korean language data can be expressed using only Korean characters without a serious problem. The enumeration method, like the Roman alphabet, is the easiest way to represent Korean characters without changing the hardware and the operating system. This method is not highly readable and would require changes in the language which may not be acceptable to users.

1. Using the Current Standard Keyboard

A program can be loaded which defines the 24 letter Korean alphabet to a character generator instead of the lower case Roman alphabet. All Korean alphabet elements and the upper case Roman alphabet characters are then available through the standard Roman character keyboard. With this method the user can use a computer in a similar manner as the users who use the Roman alphabet. In addition, well developed hardware and software can be used without critical problems. This method has been suggested by many groups of people from the time when the Korean typewriter was first developed. The only disadvantage is the breaking of traditional custom. To capitalize on developed technology and for the ease of application, more study and research should center on user acceptability of the enumeration method.

Figure 4.1 shows an example of hard copy which uses a graphic dot printer and a standard keyboard. Appendix F shows the load command program for an alternative character

```

기 840 0.1 810
LA 0.1 842 84 71 840 0.1 810 81 842 84 71 842 84 810 81
810 84 71 842 84 81 810 84 81 810 84 81 810 84 81 810
0.8 81 84.

```

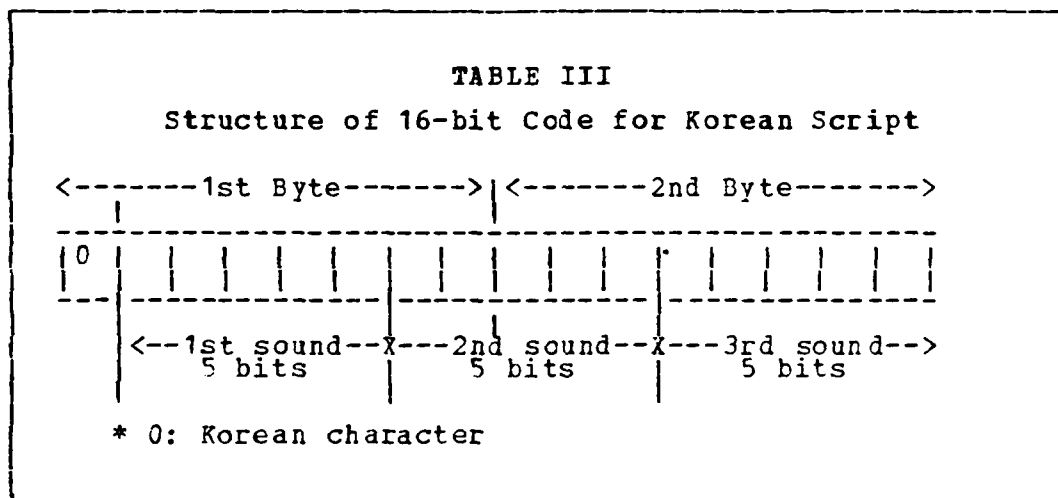
Figure 4.1 Example Using Standard Keyboard.

generator for the Korean alphabet. This program can be generated easily by the alternative character set editor, and it loads the Korean alphabet to an alternative character generator instead of the lower case of Roman alphabet.

2. Using the Capital Letters as the Initial Letter

The major difficulty with the enumeration method is poor readability. Korean users read a sentence sequentially syllable by syllable. In order to increase readability, the initial letter of each character can be written as an upper case letter to distinguish the syllable easily. Figure 4.2 shows the example using the capital letters and Appendix G represents the load command program for these letters. A special mark or altered shape of each letter also can be applied to increase a readability when an enumeration method is used.

character and Table IV explains the 16-bit code for Korean character. This code table is basically the same as the IBM 2-byte internal Korean character code [Ref. 6: p. 52]. The only difference is the arrangement. Some IBM codes represent three letters. This makes key tops (face of each key) more complex; for example, 00100 (one key top) represents L, H, and ㅈ values. The code suggested in Table IV reduces some of this complexity by limiting the possible values to no more than two for each keytop. In contrast to the example for IBM codes, the same code from Table IV represents only one value. Appendix H represents the IBM 2-byte internal



Hangul code for the Korean character.

The suggested code has several advantages. First, it is easy to sort the character order by its value since the value of each letter is in the order of the Korean alphabet. Second, it can reduce the memory space for data by using 2 bytes instead of 3 bytes for one character. Third, it is possible to edit the character directly since it does not need code conversion. Finally, since it can easily

recognize the code value of the Korean character, it helps a programmer when it is programmed.

TABLE IV
16-bit Code for Korean Script

5 bit code	1st sound letter	2nd sound letter	3rd sound letter	5 bit code	1st sound letter	2nd sound letter	3rd sound letter
00000				10000		ㅏ	ㄹㅏ
00001	ㄱ		ㄱ	10001	ㅓ		ㅓ
00010	ㄱㄱ	ㅏ	ㄱㄱ	10010		ㅑ	ㅓㅑ
00011		ㅑ	ㄱㅑ	10011	ㅕ		ㅕ
00100	ㄴ		ㄴ	10100	ㅖ	ㅏ	
00101		ㅑ	ㄴㅑ	10101		ㅓ	ㅖㅓ
00110		ㅑ	ㄴㅑ	10110	ㅗ		ㅗ
00111	ㄷ		ㄷ	10111	ㅗ	ㅓ	ㅗ
01000	ㄷㄷ	ㅏ		11000	ㅛ		ㅛ
01001	ㄷ		ㄷ	11001	ㅜ		ㅜ
01010		ㅓ	ㄷㅓ	11010	ㅜㅜ	ㅓ	
01011		ㅓ	ㄷㅓ	11011	ㅜ	ㅓ	ㅜ
01100		ㅓ	ㄷㅓ	11100	ㅟ	ㅏ	ㅟ
01101		ㅓ	ㄷㅓ	11101	ㅟ	ㅓ	ㅟ
01110		ㅓ	ㄷㅓ	11110	ㅟ	ㅓ	ㅟ
01111		ㅓ	ㄷㅓ	11111	ㅟ		ㅟ

* Blank: Not used

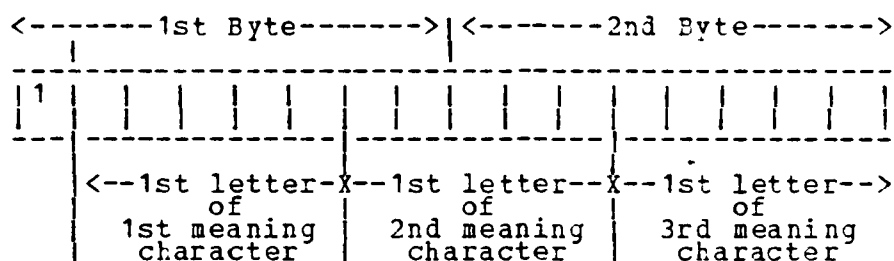
2. 16-bit Code for Chinese Characters

There is no limitation in the number of usable Chinese characters, but statistics show that 1,800-2,000 characters cover 98-99.8 percent of those which appear in newspapers and journals (Table II). Currently there are only two ways to represent Chinese characters in Korea. One method is comprised of two steps. The first step is to display all Chinese characters (synonym) which have the same sound after entering the desired sound, and the second step is to enter the Chinese character which is needed by the user via an index number matched to that character after selecting it in the display. The other method is to convert a Korean character to a Chinese character after typing a Korean character as a unit of a word, which consists of two or three characters.

The former is inconvenient and takes a long time to edit. The latter has no flexibility in that it is limited by the programmed word codes. To solve the above problem and simplify the identification of each character using a limited number of keystrokes, a 16-bit code for Chinese characters can be applied. Table V represents the structure of 16-bit code for Chinese characters.

Chinese characters represent both meaning and phonetics to Koreans. To simplify the code, all the complete meaning and sound of the Chinese characters are not needed. The Chinese characters are composed of from one to five syllables for meaning and one character for the sound. Simplicity can be achieved by employing abbreviations or acronyms for each part (meaning and sound). For example, a Chinese character (天) has a meaning as "Hea-Ven" and sound as cheon. In this case we use H of Hea, V of Ven, and C of cheon as a code for (天).

TABLE 7
Structure of 16-bit Code for Chinese Characters



* 1: Chinese character

But this method may result in duplicate codes for different Chinese characters which mean another character and may have the same value as HVC. In order to eliminate the duplicate code and to use the 3 letter code which is compatible with the 16-bit Korean character code, the following characteristics of the sound and meaning of Chinese characters are relevant: First, only 428 syllables are used to represent the sounds for all Chinese characters. That is, one sound can represent 1 to 60 Chinese characters. Second, the frequency of Korean characters used for the meaning and sound is irregular in distribution. More specific, 20% of Korean characters are used to represent the sound and meaning of 95% of Chinese characters [Ref. 11].

As a result of analyzing the 1,800 sound characters and 1,438 meaning characters used to represent the Chinese characters, Table VI and Table VIII are derived. Table VI represents the number of Chinese characters which have the same first sound letter and the same second sound letter. For example, 266 Chinese characters have the first sound letter (l), 44 Chinese characters have (l) as a first

TABLE VI
Characters Having Same 1st & 2nd Letter Sound

2nd letter		1st letter												
		ㄱ	ㄴ	ㄷ	ㄹ	ㅁ	ㅂ	ㅅ	ㅇ	ㅈ	ㅊ	ㅋ	ㅌ	ㅍ
ㅏ	ㅏ	44	8	26	17	20	27	62	25	41	22	17	11	27
	ㅑ	12	6	8	4	13	12	4	5	11	7	8	2	9
	ㅓ				9				17					5
	ㅕ													
ㅓ	ㅓ	18		1			11	42	12	59	28			7
	ㅗ	1						6		14	2			
	ㅛ	34	4		26	12	12		43				7	18
	ㅜ	15							4				6	3
ㅗ	ㅗ	34	5	36	13	18	21	28	16	28	15	9	10	24
	ㅛ	22							7	4				24
	ㅜ	1						2						
	ㅠ	4	2		2			1	2	1	2	1		8
	ㅡ	8			3	5			13				4	3
ㅜ	ㅜ	32		6	5	16	33	44	22	22	18	5	4	6
	ㅠ	6							13					
	ㅑ													1
	ㅓ	3							14		6			2
ㅡ	ㅡ	5			12				23					4
	ㅑ	20	1	4	1			10	12	8	4	1		3
ㅣ	ㅣ								10					6
	ㅣ	16	1		13	14	19	37	29	36	18		9	
Subtotal		263	27	81	105	98	135	237	265	224	122	41	53	148

Total 1,800

sound letter and (___) as a 2nd sound letter. One must rearrange the sound acronym to the 5 bit code since the distribution of sound characters is irregular. Table VII depicts the rearranged code value for acronym of the Chinese sound character.

TABLE VII
5-bit Code for the Acronym of Sound Character

code	sound acronym	code	sound acronym
00000	가	10000	스
00001	거	10001	아
00010	고	10010	어
00011	구	10011	오
00100	구	10100	우
00101	나	10101	으
00110	다	10110	자
00111	라	10111	저
01000	로	11000	조
01001	마	11001	주
01010	바	11010	즈
01011	방	11011	차
01100	사	11100	타
01101	서	11101	파
01110	소	11110	하
01111	수	11111	후

In Table VII the second letter (ㄱ, ㅋ, ㆁ, ㆁ, ㆁ) describes all the group letters. For example, (ㄱ)

represents ㅏ, ㅑ, ㅓ, ㅕ, and (ㅗ); ㅓ, ㅕ, ㅗ, ㅑ, and (ㅓ); ㅓ, ㅕ, ㅗ, ㅑ, ㅓ, etc. Also (ㅗ) represents ㅏ, ㅑ group which assembles the first consonants at the left of the vowel. (ㅗ) describes ㅓ, ㅕ, and ㅗ group which assembles the first consonants above the vowel.

Since the frequencies of Korean character syllables representing Chinese character's sound and meaning are different, the frequency of Korean characters to represent Chinese characters meaning is needed to be analyzed. After analyzing the sampled 1,438 characters which are the first and the second meaning characters, Table VIII is derived which shows the number of times for a meaning character or a group to be used.

The meaning acronym value to a 5-bit code from the basis of Table VIII can be reassigned. Table IX shows the reassigned 5-bit codes representing the acronym of the meaning character. The same theory can be applied as in Table VII when Table IX is derived. As the acronym code is rearranged, the proportion of the duplicate codes can be reduced. As a result of applying these rearranged codes, Table X can be produced which shows the proportion of the duplicate codes. The pure acronym code (Table IX) represents the acronym of a meaning and a sound character as a first letter code of Korean characters (ㄱ ㄴ ㄷ ㄹ ㅁ ㅂ ㅅ ㅇ ㅈ ㅊ ㅋ ㅌ ㅍ ㅎ; 19 possible consonants), the arranged sound character acronym code (Table VII), and the arranged sound and meaning character acronym code (Table VIII).

The reasons why some duplicate codes cannot be eliminated are: First, some Chinese characters have similar meaning and sound which generates the same acronym code (22 among 1,800); and second, there are initially some Chinese characters which have the same meaning sound (12 among

TABLE VIII
Frequency of Meaning Character in Chinese Characters

1st lt	2nd letter														TOT
	┐ group		┌ group			└ group			┘ group			- group			
	┐	┌	┌	┐	┘	└	┘	┐	┘	┐	┘	-			
ㄱ	42	6	25	9	22	25	6	6	27	3	9	26	30	236	
ㄴ	1	4	2			2		1	4	1		4	4	23	
ㄷ	23	8	7		3	9		1	7			13	2	73	
ㄹ	23	9	12			24	2	1	17		1	11	3	103	
ㅁ	1	1	3	2					2		2	1	1	13	
ㅂ	10	6	5	3	7	10		2	4		1	28	31	107	
ㅅ	39	13	10	1	3	17		2	42			1	4	93	
ㅇ	24	5	11		13	11			15				10	89	
ㅈ	2	2	1		2				1			1		9	
ㅊ	32	11	21	8		15		4	17		1	11	19	107	
ㅋ	5		1			15	1	2				4		10	
ㆁ	28	6	29	4	19	15	1	2	22	4	2	57	55	246	
ㄷ	18	5	8	2		6			12			3	39	93	
ㅈ	2					1			1				4	8	
ㅊ	5	4	6			3			4			5	13	40	
ㅋ	2	1				1					1	9		14	
ㆁ	8	3	4			3			2			2	1	23	
ㅂ	4	2	1		1	1			5			1	2	17	
ㅇ	32	2	2	2		6	7					2	10	69	
	301	87	149	31	70	143	16	19	182	8	17	187	228	1438	

* 2nd letters having 10 are omitted for simplicity

TABLE IX
5-bit Code for Acronym of Meaning Character

code	1st & 2nd meaning acronym	code	1st & 2nd meaning acronym
00000		10000	ㅅ.
00001	가	10001	ㅇ
00010	거	10010	아
00011	고	10011	어
00100	구	10100	오
00101	그	10101	우
00110	ㄴ.	10110	으
00111	ㄴ	10111	ㅈ.
01000	ㄷ.	11000	중
01001	ㄷ	11001	ㅊ
01010	ㄹ.	11010	ㅋ
01011	ㄹ	11011	ㅌ
01100	ㅁ.	11100	ㅍ
01101	ㅁ	11101	ㅎ.
01110	ㅂ.	11110	ㅎ
01111	ㅂ	11111	

1,800). For these characters which have duplicate codes the users must apply the exception rule. Alternatively, the meaning of the character can be redefined to a synonym with a different acronym. For example, if the wrong homonymous Chinese character is displayed, the input operator may select another form of the homonym by keying in the full sound syllables instead of the acronym.

TABLE X
Proportion of Duplicate Code

Number of character	pure acronym code	Rearranged sound character code	Rearranged sound and meaning character code
1,800	7.5%	2.3%	1.8%

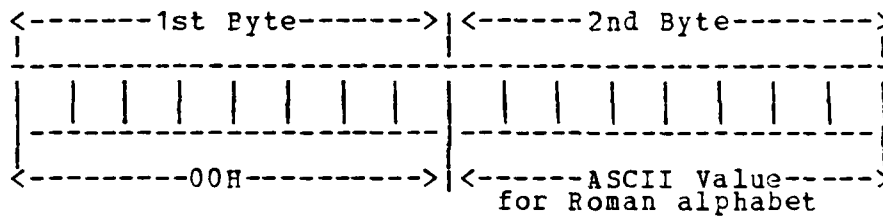
3. 16-bit Code for Roman Alphabet and Symbols

In order to use the three mixed kinds of a character code, and simplify the I/O controller, and unify the word, 16-bit codes for the Roman alphabet and symbols should be generated by only one keystroke. For data communication and for familiarity, adding only the default byte (00H) to ASCII code, 16-bit code for Roman alphabet, symbols, and control characters can be defined. When one uses only Roman alphanumeric characters, one can easily convert this 16-bit code to ASCII code. Table XI shows the 16-bit code for Roman alphabet and symbols.

4. Keyboard for 16-bit Code

As it is mentioned in the previous chapter, the biggest issue is how to enter all Chinese characters, Korean characters, and the Roman alphabets with a simple keyboard. In order to implement the 16-bit code to keyboard, one would have to make the keystrokes which generates "1" or "0" as a Chinese character function key (bit 1), three 5-bit codes (00000-11111) for Chinese and Korean characters (bits 2-16) and 16-bit code for Roman alphabet and symbols. In this case 33 more keys than the common Roman alphabet keyboard are

TABLE XI
Structure of 16-bit Code for Roman Alphabet



needed. To reduce the number of keys, one more function key can be added for the Roman alphabet which generates 16-bit code as a 5-bit code key. However, the user identification will be complex because there must be 4 or 5 letters on each key top. Table XII explains the 4 alternatives. Alternative I includes 32 Roman alphabet on 5-bit key tops using a Roman alphabet function key and Alternative II excludes Roman alphabet on 5-bit key top. Alternative A uses the acronym of sound and meaning characters and Alternative B uses only the acronym of sound character for the sound and meaning characters.

To select the best one, the authors can use the following criteria: flexibility of hardware and software design, hardware efficiency, ease of maintenance, system reliability, user characteristics, number of keystrokes, number of duplicate codes, and complexity of recognizing a certain keystroke [Ref. 2]. In the opinion of authors' alternative II-B should be selected for the ease of explanation and understanding. By selecting Alternative II-B, a user can type the three characters simultaneously. For example, to type "School is 학교 (Hak-kyo) in Korean and 學校 (Hak-kyo) in Chinese character", one can type directly

TABLE XII
Leveled Letters on 32 Key Tops

Alternative I		Alternative II	
I-A	I-B	II-A	II-B
<div style="border: 1px dashed black; padding: 5px; display: inline-block;"> 1 2 3 4 5 6 </div>	<div style="border: 1px dashed black; padding: 5px; display: inline-block;"> 1 2 3 4 5 </div>	<div style="border: 1px dashed black; padding: 5px; display: inline-block;"> 2 3 4 5 6 </div>	<div style="border: 1px dashed black; padding: 5px; display: inline-block;"> 2 3 4 5 </div>

Legend: 1 : Roman alphabet
 2 : First letter of Korean character
 3 : Second letter of Korean character
 4 : Third letter of Korean character
 5 : Acronym of Chinese sound character
 6 : Acronym of Chinese meaning character

only Roman alphabets in the previous sentence without a function key. Two syllables of "school" are formed from: In Korean, the first syllable is "ㅏ" selected from the position of first sound letter, "ㅑ" from the second sound position, "ㄴ" from the third position. The second syllable is "ㄴ, ㄷ, Default". Then, in Chinese, press the Chinese character function key which generates "1" as the first bit. The first syllable is "ㅏ" which is the acronym of first meaning character and "ㅑ" which is the acronym of second meaning character and "ㅏ" which is the acronym of sound character. The second syllable is "ㅏ ㅑ ㅏ". After typing Chinese characters, user must release the function key to type Korean characters. Table XIII explains the above example.

As a result of the above example, the computer generates the following codes in hexadecimal: 0053(S),

TABLE XIII

Typing Procedures for Mixed Characters

To type "school, 학 교, 學校", the following procedures should be followed:

1. Type "school" by one keystroke for each character without a function key.
2. In Korean, to type "학 교",
first, type "ㅎ, ㅌ, ㄱ",
second, type "ㄱ, ㅍ, Default".
3. In Chinese, to type "學 校", at first, press the function key, then type "ㅂ, 우, ㅎ",
"ㅎ, ㄱ, ㄱ" (Table IX) since for "學",
the meaning is "배 울" and the sound is "학", and
for "校", the meaning is "학 교" and the sound
is "교".

0063(c), 0068(h), 006F(o), 006F(o), 006C(l), and 0 (Korean character), 11111(ㅎ), 00010(ㅌ), 00001(ㄱ); that is 0111 1100 0100 0001(7C 41:학), 0 (Korean character), 00001(ㄱ), 10010(ㅍ), 00000(Default); that is 0000 0110 0100 0000(06 40:교), and 1 (Chinese character), 01010(ㅂ), 10100(우), 11110(ㅎ); that is 1010 1010 1001 1110(AA 9E:學), and 1 (Chinese character), 11110(ㅎ), 00010(ㄱ), 11110(ㄱ); that is 1111 1000 0101 1110(F8 5E:校).

5. Operating System for Input and Output

To apply the suggested system, it is needed to redesign the operating system for input and output control.

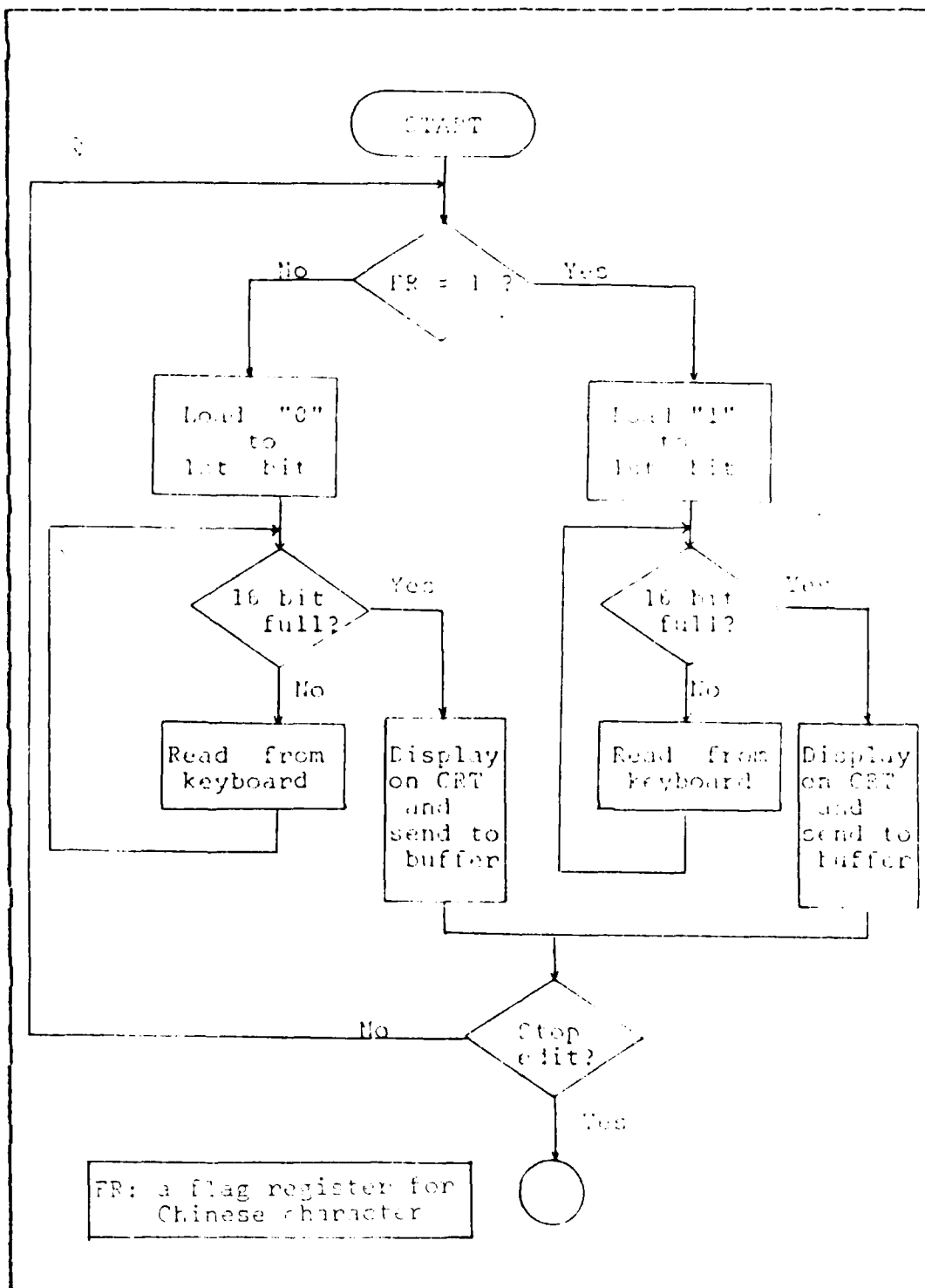


Figure 4.3 Flowchart of Input and Output Controller.

Figure 4.3 shows the flowchart of input and output control. First, the input and output controller has to distinguish whether the Chinese character function key is "0" or "1". A flag register can be used for Chinese character function key. For example, if the flag is "1", then "1" is loaded to the first of the 16 bit register, and multiple 5-bit codes are read until the register is full. If it is full, a character is displayed on the CRT, and the 16-bit code is sent to a buffer as data. Otherwise the flag is "0", and then "0" is loaded into the register, and 5-bit codes are read for a Korean character. A 16-bit code is used for a Roman alphabet character or a symbol code until the 16-bit register is full. If the 16-bit register is full, the identified character is displayed and sent to a buffer as data. If the "stop edit?" condition shown in Fig 4.3 is "no", the input and output controller makes a loop to read a code, displays a character and sends a character code to a buffer.

This system will make the use of Korean language commands and programs easier to use than those presently available. To achieve the above goals, a compiler and interpreter, as well as the operating system will require redesigning. This system will require the complete rewriting of all software currently used. The economic impact of this on the Korean people will be enormous.

6. Design Considerations for Character Generation

There are two shapes of characters used in Korea: Gothic (Figure 4.4) and Brush type (i.e., Ming style: Figure 4.5) [Ref. 7: p. 34]. To generate the above shapes of characters, several methods of a character generation can be considered. To select the best method for Korean and Chinese characters, one can use the following five criteria: speed, space, quality, flexibility, and cost. Speed is a double standard: speed of creation may range from a few minutes to

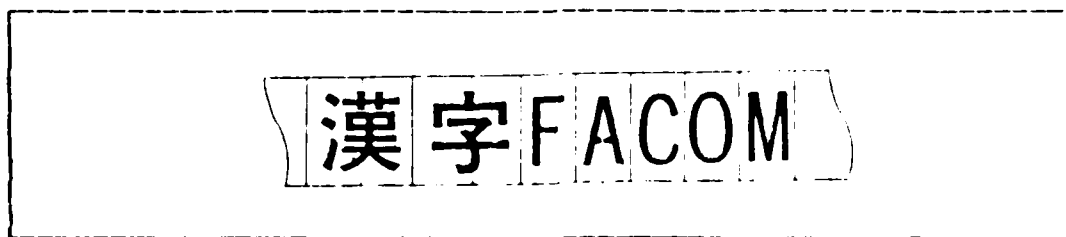


Figure 4.4 An Example of Gothic Type.

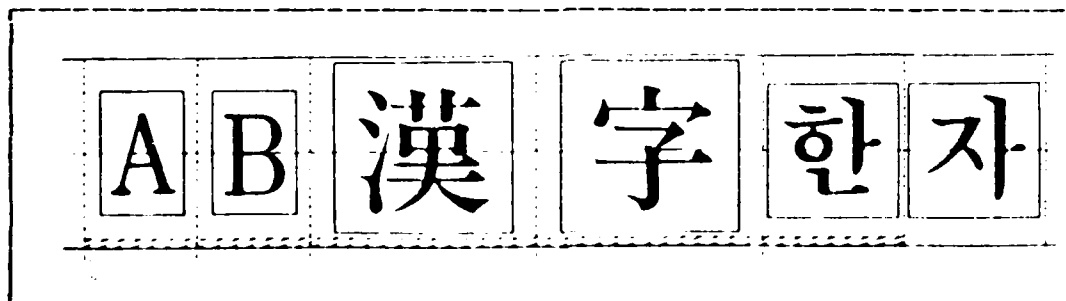


Figure 4.5 An Example of Brush Type.

a few hours, while speed of production should go beyond 1000 characters/sec depending on type size and device resolution. Space refers to the average size of the code for one character as well as the size of the internal buffers often needed for decoding. Quality is proportional to the largest dot matrix which can be used to decode a character; it should not be confused with the resolution of the output device. For a given type size, the resolution sets the definition, that is the size of the matrix to be used; definition, hence type size, is bounded by the quality of the code. Flexibility refers to the different automatic modifications which are supported by the code: scaling, rotating, family variations (as going from light to bold). Cost is selfexplanatory [Ref. 12: p. 240].

Obviously the five criteria above are not independent. Figure 4.6 shows the interrelationship of criteria in designing a character generator [Ref. 12: p. 241]. The most desirable feature is indicated by the direction of the arrow. Solid (resp. dotted) lines indicate agreement (resp.

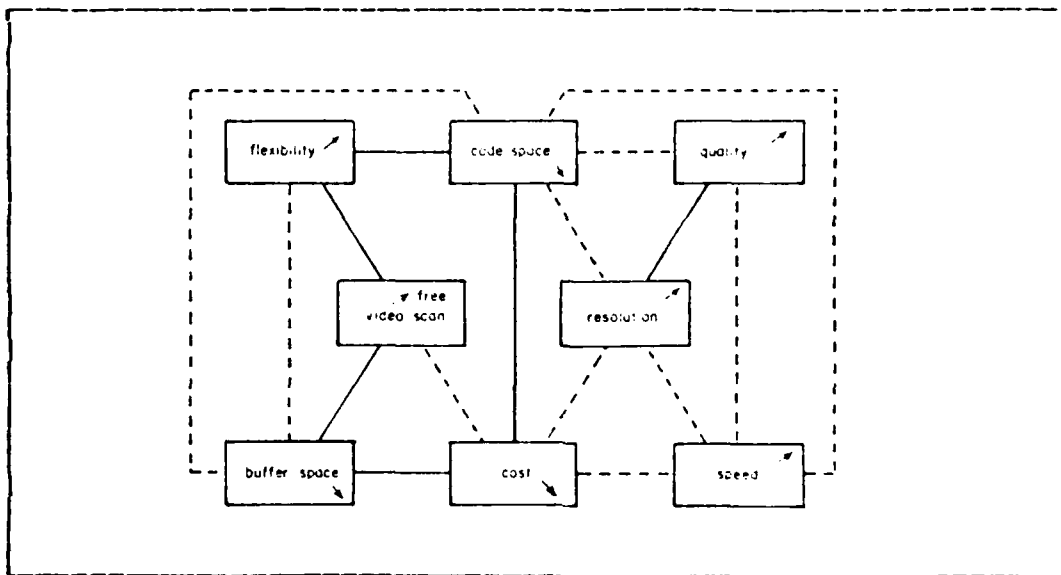


Figure 4.6 Criteria in Designing a Character Generator.

contrariety) between the variation of the factors. The design of a digital character generator is an engineer's task whose goal is to strike the appropriate balance between the specifications for those five criteria combined with the characteristics of the production device, resolution and scanning, and the necessity of operating the corresponding creation station.

Table XIV [Ref. 12: p. 268] gives a summary of the main characteristics of the coding methods that the engineer can utilize [Ref. 12: p. 269]. As the characteristics of Korean and Chinese characters are compared to Table XIV, it should be apparent that the bit map method is the most

TABLE XIV
Comparative Table for Performance

Method	Code space (bits)	Buffer space (bits)	Flex- ibility	Video scan	Reso- lution (n)	Qty	Spd
Bit map	$\frac{2}{n}$	0	-	+	≤ 50	-	+++
Run- length	$k \cdot n \cdot \log_2 n$	n	-	+	≤ 100	-	++
Chain- link	$6 \cdot k \cdot n$	$\frac{2}{n}$	-	-	≤ 100	-	+
Diffe- renti- al run length	$\frac{6 \cdot k \cdot n + b(\log_2 n + c)}{2}$	n	-	+	≤ 100	-	+
Spline	$k \cdot \log_2 n$	$k \cdot \log_2 \frac{n}{m}$	-	-	-	-	-
Struc- tural	$k' \cdot \log_2 n$	$\frac{2}{n}$	-	-	-	-	-

* Qty: Quality
* Spd: Speed

Legend: +: Good
-: Bad
b: The number of birth point in the character
c: A constant taking care of bookkeeping
n: The size of matrix
k: The average number of runs per matrix
line or column (a number of the simplicity
of the character shapes: approximately 4
for Roman body-text fonts, higher than
10 for Chinese characters)

appropriate one for this application. It reduces code space and buffer space. It has good video scan, high speed, and highly readable low quality printing. Unfortunately this method lacks flexibility. However for all the other aforementioned reasons, the bit map method is commonly used for Korean and Chinese characters.

Because of the complexity of Korean and Chinese characters, at least a 16 by 16 resolution is required for Korean characters, and a 24 by 24 resolution for Chinese characters. 32 by 32, 64 by 64, 80 by 80, 96 by 96, and 128 by 128 resolutions are desirable when much more beauty is required and also when larger character sizes are to be produced. However, if these characters are displayed on a CRT, with 32 by 32 resolution, with the size of each character 7-10 mm square, this should be sufficient.

It is the authors' opinion that the less expensive 32 by 32 resolution CRT should be used for softcopy. The reason for this is that the price of the memory component required to hold the character definitions is continually getting less expensive. However stronger motivation is that high-speed and flexibility of typing is then possible [Ref. 1: p. 828]. IBM corporation uses 16 by 16 resolution for Gothic and 24 by 24 for Brush type Korean character syllables [Ref. 8: p. 2]. FACOM corporation uses 30 by 30 dots for Korean and Chinese characters and 24 by 30 dots for Roman alphabet, symbols, and Korean alphabet (letters) on a laser printer [Ref. 7: p. 47]. As cheaper dot matrix and laser printers find their way into the marketplace, the quality of characters will become less of an issue. Presently there are few problems with the quality of representing Korean characters that cannot be solved through the additional expenditure of money. For the definition of each character, the authors have presented two alternatives; software and hardware (character generator). In order to increase speed and usability, a hardware-oriented character generator is best. If cost and flexibility are the criterion, software-oriented character definition programs are better.

7. Memory Space for Character Definition

To represent Korean and Chinese characters, one needs to code 5,000 characters for a character definition: (2,400 Korean characters; 1,800 Chinese characters; 800 user definable characters, Roman alphabet, or symbols). If one uses a 24 bit by 24 bit matrix font size for each character, at least 360 K bytes are required ($3 \text{ byte} * 24 * 5,000$) for character definition and 128 K bytes are required (64 K for 16 bit address * 2 byte for an address of character definition memory) for a look-up table. The total memory requirement is 488 K bytes.

A large memory space is required for the definition of the characters. Data compression of these characters can be considered for two different purposes: data transmission, and computer storage and output. Here one is mainly interested in the latter case, where the main point is the total data amount to be stored. The method of data compression of Chinese characters can be classified by using the methodology listed in Table XV [Ref. 1: p. 820].

There is a problem associated with the enlargement and alignment of character patterns. The clarity of a character depends on the size of the reproduction. If a large size is required the resolution must be high. Otherwise, stepwise zigzags appear which to some people are unbearable. Therefore, all the patterns of different font sizes must be stored. This is uneconomical. Reproducing different character sizes from the same data is desired. However, the enlargement and shrinking of character patterns from a single set of data is quite difficult, because, if the addition or the deletion of a bit by the interpolation is not done properly, it has a negative influence on the aesthetics. In enlargement, the smoothness of an edge is particularly important, while in shrinking the gap between strokes must be carefully maintained.

TABLE XV
Varieties of Data Compression Methods

transmission	{	page unit	{	run-length coding
		character unit	{	two-dimensional predictive coding
				coding by scan line pattern unit
memory & reconstruction	{	dot pattern representation	{	coding by m by n block pattern unit
		stroke representation	{	checker board sampling
		contour coordinate coding	{	hexagonal board sampling
		contour following coding		
		mathematical equation for strokes		
enlarge/shrink	{	synthesis from partial character patterns		

A comparative review of the options contained in Table XV with regard to determining memory size is very difficult. This is because the requirements for character print qualities are quite different depending on each method. Simplicity in the hardware and software implementation of the compression and reconstruction of characters is a very important consideration. Generally speaking high data compression methods need complex hardware and longer times for reconstruction. Therefore, the tradeoff to be considered is between the data compression ratio and the memory size. This represents the classic economic tradeoff between the hardware/software cost with regard to the speed of character regeneration.

Because the price of the memory component is becoming less expensive the high-speed simple reconstruction method is preferred despite the necessarily large size memory. Many commercial machines have adopted this concept, and store the character dot patterns as they are without any data compression. For example, IBM machines use only a 12 by 24 font size for simple letters (Roman and Korean alphabet and symbols) instead of a 24 by 24 font [Ref. 8: p. 17]. The FACOM machines use the software definition of the second level of Korean and Chinese characters which are not used frequently for data compression [Ref. 7: p. 12]. Because of the reduction in the price of memory, the marketplace has shifted towards providing direct character storage, i.e. a large memory, instead of utilizing data compression.

V. EVALUATION OF SUGGESTED METHODS

The principal problems in current editing technology for Chinese and Korean characters were detailed in Chapter III. Fundamentally, the problems cause user inconvenience, require lengthy input procedures, and result in complex update requirements. The authors' suggested methods will solve most problems which are encountered in Korean language data processing. More research and development remain in the following areas:

First, in an enumeration method, there is no problem except low readability to Koreans and the inability to represent Chinese characters. In this case, Chinese characters are ignored because Korean language data can be represented through the use of only Korean characters without serious problems. Low readability is caused by unfamiliarity and the unbalanced shape of each letter when written by an enumeration method. With a minor change of shape of the letters, this method will eliminate the above problems.

Second, the 16-bit code for the three kinds of characters requires the consideration of the following problems:

1. The 32 key tops are complex since each key top represents three or four letters and acronyms. One solution to this problem is to use lighted, changeable key tops which represent only one letter or acronym at a moment according to the function keys and the order of keystrokes(i.e., 1st, 2nd, 3rd letter and acronym).
2. The user must remember whether a letter to be typed is the first, second, or third letter, and whether it is an acronym of a sound or a meaning character.

3. If a user does not know the meaning of a certain Chinese character to be typed, one must look up a table which shows all meanings and sounds of all possible Chinese characters.
4. In typing Korean characters which consist of only first and second letters, a user has to hit the default key to make a 16-bit code. Instead of second letter and default keys, one can use twenty one more second letter keys which generate 10-bit code as the second and third letters. Unfortunately, this will make the keyboard more complex.
5. Regardless of the authors' analysis of sound and meaning characters and careful rearrangement of these codes, duplicate codes still exist. This is because of the irregularities caused by a natural evolution of sound and meaning characters for over 2,000 years. Generally 3,000 or more Chinese characters will cause duplicate codes to increase proportionally. In order to eliminate the duplicate codes, the Korean language committee needs to take measures to clarify the meanings of the Chinese characters that cause duplicate codes to exist.

Before the actual construction of the suggested system, an economic (Cost/Benefit) analysis needs to be considered. Given the r % discount rate and the various yearly costs and benefits estimated by past data, Table XVI [Ref. 14] shows the following formula which can be used to derive the net present value of this project: This simply states that the net present value (NPV) is equal to the sum of the differences between benefits (B) and costs (C) in each year (i) of the project life (T), divided by the relevant factor (r) for that year. The current estimate of the market size for word processing in Korea is \$ 2.5 million annually (Korean Daily Times, Sep 10 1984). But this estimate will be in inverse

TABLE XVI
Net Present Value Formula

$$NPV = \sum_{i=1}^T \frac{(B_i - C_i)}{(r)^i}$$

Legend: NPV: Net Present Value
B: Benefit
C: Cost
i: Each year
T: Project life
r: Relavent factor

proportion to the price of the system and will be in direct proportion to the usefulness and the user-friendliness until maturation.

In the above formula, the market price of the system influences the benefits for manufacturers and costs for users. This system is feasible when the net present values are positive for both manufacturers and users. If the benefits for manufacturers and the costs for users are constant in a system, the main problems will be:

1. How to minimize the costs for manufacturers
2. How to maximize the benefits for users

To solve the above problems, the best approach will be to make an efficient and user friendly system for Korean language data processing. This will increase the number (N) of systems sold, and increase the individual productivity of the users.

There are many factors and constraints which cause high cost in implementing this method. Among these, the following

three factors affect the cost performance ratio for both manufacturers and users:

1. The initial design cost: For this system, an organization has to invest initially for the design of about 5,000 Korean and Chinese character patterns, and the system software and hardware. As the number of systems produced by a manufacturer is increased, the unit cost of each system will be decreased as the costs are spread over more units.
2. Cost for character generator: As mentioned in Chapter IV, one needs about 500 K bytes memory capacity for these character definitions. The cost of memory is decreasing and the speed is increasing as technology is being developed. This cost is an initial cost to users when buying a system.
3. Cost for hardcopy: One can consider three kinds of printer for hardcopy: dot matrix printer, chain printer, and laser printer. It is not practical to use a chain printer for our system since the chain will be approximately twenty meters long (5,000 character * 4 mm per each character) and it would be prohibitively slow. Currently, dot matrix printers and laser printers cost more than chain printers, but they are the only viable option.

Among the three kinds of cost, the third one is the most serious since the cost of hardcopy is increasing as its use increases. Recently laser printers have become more popular for these characters because of the good quality, high speed and decreasing price. Comparatively though, laser printers are still relatively expensive.

VI. RECOMMENDATION AND CONCLUSION

As the demand for data processing in Korea increases, users will continue to encounter more and more problems in utilizing the Korean language for data processing. The current methods of convergence, display and select to implement the Korean language in data processing must only be considered as interim measures due to their inefficient and time consuming means of data entry. In order to prevent this problem from becoming more complicated due to the development of various new implementations forwarded by independent research, a standardized system must be developed.

This study examined two possible solutions for using the Korean language in data processing. The enumeration method is technologically feasible, inexpensive, easy to implement, but could not be used for applications within the Korean data processing environment. This is because it results in a textual form of Hangul that is unfamiliar to most Korean people. Therefore the current enumeration method is not a feasible solution to the Korean data processing problem.

The second method examined was based on a 16-bit code representation of Korean, Chinese characters, and the Roman alphabet. This method was found to possess all the advantages currently realized by the EBCDIC or ASCII code representation of western countries. The only drawback to this system is that it might not be cost effective based on current technology. However, due to the rapid development of hardware and software technology, a cost effective means should be available within the next few years.

In order to accelerate the determination of a thorough broad based solution to the Korean data processing problem, the Korean government must organize and charge a national

level committee with the responsibility for investigating the problem and determining a viable solution. This study with its proposal of a 16-bit character code should be provided to that committee for further examination. This proposal represents a concept that could eventually lead to a long term viable solution to the data entry and processing problems of using the Korean language.

THE EVOLUTION OF CHINESE CHARACTEES

Changes
after
Han Dynasty
(after
A.D. 588)

51

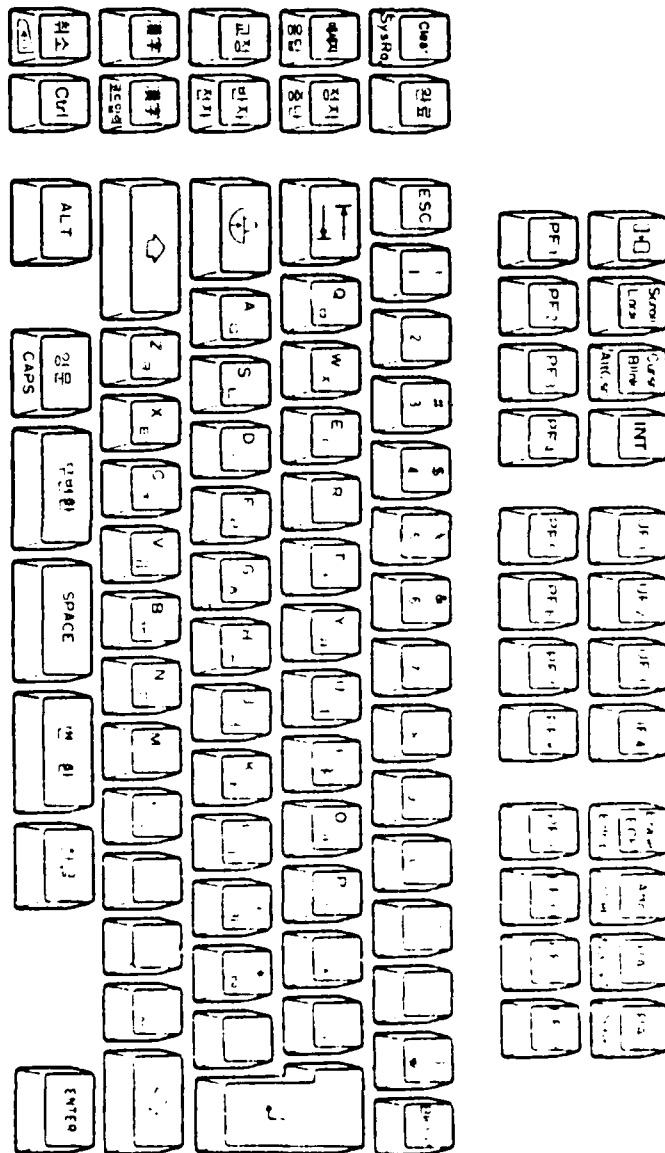
APPENDIX B **EBCDIC INPUT CODE**

		00				01				10				11				
Hex 1		00	01	10	11	00	01	10	11	00	01	10	11	00	01	10	11	Hex 0
Bits 4 5 6 7		0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	
0000	0	NUL	DLE			SP	X									S	0	
0001	1	SOH	SBA											A	J		1	
0010	2	STX	EUA		SYN									B	K	S	2	
0011	3	ETX	LC											C	L	T	3	
0100	4	VCS	ENP	INP										D	M	U	4	
0101	5	PT HT	NL	LF	TRN									E	N	V	5	
0110	6		BS	ETB										F	O	W	6	
0111	7			ESC	EOT									G	P	X	7	
1000	8													H	Q	Y	8	
1001	9		FM		LDEL									I	R	Z	9	
1010	A				LINE													
1011	B	VT	SFE	FMT														
1100	C	FF	DUP		RA													
1101	D	CR	SF	ENQ	NAK													
1110	E		FM IRS															
1111	F		ITB	BEL	SUB													

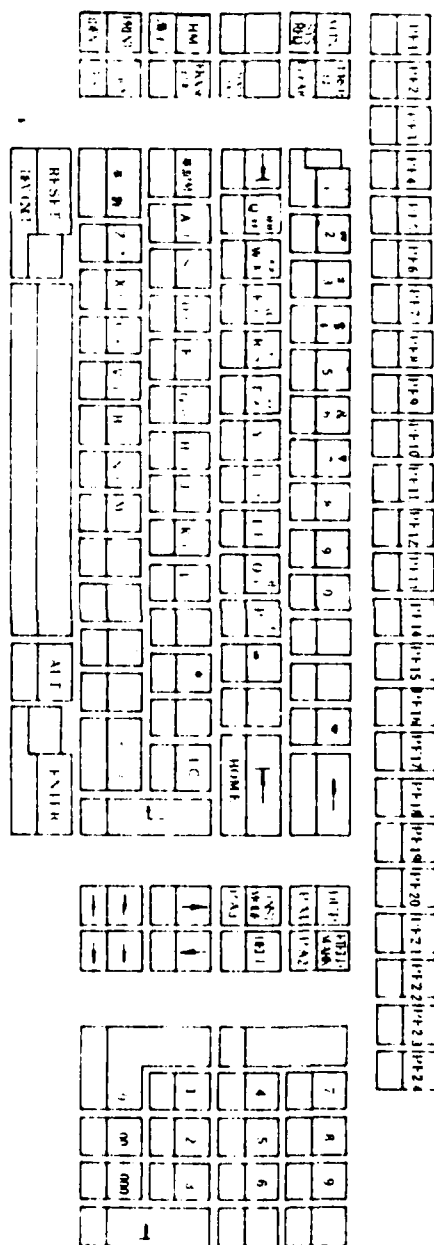
APPENDIX C
MDS INPUT CODE

[illegible]

APPENDIX D
IBM MULTISTATION 5550 KEYBOARD



APPENDIX E FACCM OS IV (KEF) KEYBOARD



APPENDIX F

LOAD COMMAND PROGRAM FOR CURRENT KEYBOARD

```

1(a"KC")
1" "00;
1"A"00081422417F4141;
1"D"007C22212121227C;
1"E"007F40407C40407F;
1"F"007F40407C404040;
1"G"001E21404047211E;
1"H"004141417F414141;
1"I"003E08080808083E;
1"J"003E41414141413E;
1"P"007E41417E404040;
1"Q"0C3E41414145423D;
1"R"007E41417E444241;
1"S"003E41403E0.413E;
1"T"007F080808080808;
1"U"004141414141413E;
1"W"0041414149495522;
1"Y"0041221408080808;
1"a"00007F414141417F;
1"d"00001C224141221C;
1"e"00007F404040407F;
1"f"00C07F01017F407F;
1"q"0000087F003E413E;
1"u"000008080808087F;
1"i"0000203E20203E20;
1"o"0C004242427E4242;
1"p"0000090909790909;
1"u"000041417F41417F;
1"r"00007F0101010102;
1"s"0C0040404040407F;
1"t"0000080808142241;
1"u"0000011F01011F01;
1"u"00007F0808142241;
1"v"000022222222227F;
1"! "0008080908080008;
1"1"0C0818280808083E;
1"<"0001061860180601;
1"2"0C7E21213E21217E;
1"C"001E21404040211E;
1"J"000E040404044438;
1"K"00434C5850584C43;
1"L"0C4040404040407F;
1"M"0041535549414141;
1"N"0041615149454341;
1"V"0041412222141C08;

1"1"0000080808080808;
1"m"000000000000007F;
1"n"00007F0808080808;
1"v"00007F222222227F;
1"x"00007F007F40407F;
1"z"00007F011F010204;
1" ,"0000000000302040;
1"."00000000001818;
1">"002018060106182C;
1"'"000022147F1422;
1"'"0000494977494977;
1"'"000008087F0808;
1"'"0000097909097909;
1"["00007F2222225549;
1"]"0C00774444444477;
1"("0000602010;
1")"000E10107010100E;
1"#"0C24247E247E2424;
1"$"00083E483E093E08;
1"%"00515204081C2343;
1"%"00050408;
1"%"0C04020101010204;
1"%"00000000000000FF;
1"%"001C20404040201C;
1"%"0018180C001818;
1"%"0000007F;
1"%"0000222236494949;
1"0"0C1C22454951221C;
1"2"003C42010E30407F;
1"3"007F02040E01413E;
1"4"00040C14247F0404;
1"5"0C7F407E4101413E;
1"6"001E21407E61211E;
1"7"0C7F0102040E1C2C;
1"8"003E41413E41413E;
1"9"003C42433C01423C;
1"="000C0C00000C081C;
1"? "003E410608080008;
1"@"0C242424;
1"\ "0000771111111122;
1"^"0038444438454639;
1"_ "0000007F007F;
1"~"004141FF41495522;

```

APPENDIX G

LOAD COMMAND PROGRAM FOR CAPITAL LETTER KEYBOARD

```

l(a"lee")
l" "00;
l"! "0000007F01010101020C;
l"#"000000007F0314224141;
l"$ "0000000083E483E093E08;
l"% "00000000007F0022227F;
l"^ "0000000000001818;
l"_"000000009090979090909;
l" ` "000000005052527700505;
l"/ "000000009097909790909;
l"0 "000000005701525504545;
l"1 "00FF0101010101010101;
l"2 "007F414141414141417F;
l"3 "007F0808142241414141;
l"4 "007F002424242424247F;
l"5 "00494949497F4949497F;
l"6 "0000004040407E404040;
l"7 "0000000000080808087F;
l"8 "000000000000000007F;
l"9 "000000001010101030579;
l": "000000001701121274141;
l"; "0000000007F414141417F;
l" "0000000007F404040407F;
l"b "0000000075151577454577;
l"c "0000000007F007F40407F;
l"d "0000000007F013001020C;
l"e "000000001007F08142241;
l"f "0000000077090909112244;
l"g "0000000077151575454577;
l"h "000000000771111714172;
l"i "00000000040404040407F;
l"j "000000001C007F003E413E;
l"k "00000000040814224141;
l"l "00000000121232404949;
l"m "000000000721212354949;
l"n "00000000041417F41417F;
l"o "000000001C224141221C;
l"p "000000000771015754077;
l"q "0000000007F01017F407F;
l"r "000000006404F40464975;
l"s "007F404040404040407F;
l"t "007F000007F404040407F;
l"u "007F0101013F01010101;
l"v "001C007F031422414141;
l"w "000000004444447C4444;
l"x "000000001017F09091151;
l"y "0000000007F1414141414;
l"z "000000002023E023E0202;
l" "000000001111117F0101;
l" "0000000040424277C0404;
l" "0040404040404040407F;
l" "001C007F001C2241221C;
l"s "00010204081422414141;
l"t "007F242424344A494949;
l"u "0000000000C141414147F;
l"v "0077444444444444477;
l"w "00414141417F4141417F;
l"x "001C224141414141221C;
l"y "00000040407C407C4040;
l"z "007F0101017F4040407F;
l" "00000001030501010107C;
l"(" "00000007F02040E01413E;
l") "0000000040C14247F0404;
l" "00000003C42010E30407E;
l"i "00000001C22405E51211E;
l"u "00000007F405E5101413E;
l" "000022147F1422;
l" "00181800001312;
l" "00000007F;
l" "00000008067F0808;
l"< "0001061860130601;
l"=" "000000000000000810;
l"> "0020180601061820;
l"? "003E410606060006;
l"J "003E41413E41413E;
l" "003C42433C01423C;
l"L "0038444436454639;
l" "001C22454951221C;
l"N "00080808080808008;
l"O "007F010204031020;
l"P "0051620408102343;
l" "0010204040402010;
l"\ "00060406;
l"] "0004020101010204;
l"_ "0000007F007F;
l" "00414141414141414141;

```

APPENDIX H
IBM 2-BYTE INTERNAL HANGUL CODE

Hex 134	A	B	C	D
0000	Not Used-1	Restricted-1	Restricted-1	Restricted-1
0001	Full Code-1	Restricted-1	Restricted-1	Full Code-1
0010	1	Full Code-1		
0011	2			
0100	3			
0101	4			
0110	5			
0111	6			
1000	7	Restricted-1		
1001	8	Restricted-1		
1010	9	1		
1011	A	2		
1100	B	3		
1101	C	4		
1110	D	5		
1111	E	6		
1000	F	Restricted-1		
1001	10	Restricted-1		
1010	11	1		
1011	12	2		
1100	13	3		
1101	Not Used-3	4		
1110	Not Used-3	5		
1111	Not Used-3	6		
1100	Not Used-3	Restricted-1		
1101	Not Used-3	Restricted-1		
1110	Not Used-3	1		
1111	Not Used-3	2		
1110	Not Used-2	Not Used-3		Not Used-3
1111	Restricted-1	Not Used-3		Restricted-1

LIST OF REFERENCES

1. Nagao, Makoto, "Data Compression of Chinese Character Patterns", Proceedings of the IEEE, Vol. 68, No. 7, pp. 816-819, July 1980.
2. Mackenzie, C. E., Coded Character Sets, History and Development, Addison-Wesley Publishing Company, 1980.
3. Dong-A Publishing Company, 1800 Characters of Sino-Korean Textbook, 1978.
4. Department of Computer Engineering., Seoul National University, English/Hangeul/Chinese Character Display Controller Design Using Address Conversion Technique and DMA, by C. M. Kim and H. Y. Hwang, 17 May 1982.
5. Rev. Kim, Jacob Chang-Ui, "The Study of Etymology, Classification and Signification of the Chinese Character", Sino-Korean, Herald Printers, Monterey, California 1984.
6. Lee, M. N., A Study on 2-Byte Hangeul Database Systems, M.S. Thesis, Yonsei University, Seoul, June 1982.
7. FACOM OS IV KEF Manual, 1st ed, V.1, FACOM Korea Corp., June 1982.
8. IBM Multistation Manual, 1st ed, IBM Korea Corp., 1984.
9. Grant, Bruce K., A Guide to Korean Characters, Hollym International Corp., 1979.
10. Defense Language Institute Foreign Language Center, Korean: Resource Module I Sounds and Hangeul, September 1980.
11. Hangeul Academy, A History of Hangeul Academy, December 1981.
12. Coueignoux, "Character Generation by Computer", Computer Graphics and Image Processing, pp. 240-269, 1981.
13. Foley, James D. and Van Dam, Andries, Fundamentals of Interactive Computer Graphics, Addison-Wesley Publishing Company, March 1983.

14. Fin-Dor, Phillip and Jones, Carl R., Information
Resources Management (Draft).

BIBLIOGRAPHY

Hamacher, V. Carl, Vranesic, Zvonko G., and Zaky, Safwat G., Computer Organization, McGraw-Hill Book Company, 1978.

Hall, Douglas V., Microprocessors and Digital Systems, 2nd ed., McGraw-Hill Book Company, 1983.

Madnick, Stuart E. and Donovan, John J., Operating Systems, McGraw-Hill Book Company, 1974.

Panco, Raymond R., "Kanji Keyboard Chaos", Computerworld, 28 November 1983.

People's Republic of China Library Academy, Chinese Character Code for Information Interchange, 1969.

INITIAL DISTRIBUTION LIST

	No.	Copies
1. Defense Technical Information Center Cameron Station Alexandria, Virginia 22314	2	
2. Library, Code 0142 Naval Postgraduate School Monterey, California 93943	2	
3. Department Chairman, Code 54 Department of Administrative Sciences Naval Postgraduate School Monterey, California 93943	1	
4. Department Chairman, Code 52 Hq Department of Computer Science Naval Postgraduate School Monterey, California 93943	1	
5. Professor Michael J. Zyda, code 52 Department of Computer Science Naval Postgraduate School Monterey, California 93943	1	
6. Professor Carl E. Jones, code 54 Js Department of Administrative Sciences Naval Postgraduate School Monterey, California 93943	1	
7. Library Officer Korea Military Academy Seoul, Korea 130-09	2	
8. Office of the Defense Attache Embassy of the Republic of Korea 2320 Massachusetts Avenue, Northwest Washington D.C. 20008	1	
9. Major Chong Hae Kim Planning and Management Staff Department Republic of Korean Army Headquarters Seoul, Korea 130-09	5	
10. Major Sung Woo Kc Planning and Management Staff Department Republic of Korean Army Headquarters Seoul, Korea 130-09	5	
11. Major Jang Hun Lee, SMC 1366 Naval Postgraduate School Monterey, California 93943	1	
12. Captain Hee Young Lee, SMC 2726 Naval Postgraduate School Monterey, California 93943	1	
13. Captain Hwa Dal Song, SMC 2748 Naval Postgraduate School Monterey, California 93943	1	

- | | | |
|-----|--|---|
| 14. | Captain Kwang Jun Choi, SMC 2870
Naval Postgraduate School
Monterey, California 93943 | 1 |
| 15. | Captain Jai Eun Jang, SMC 1032
Naval Postgraduate School
Monterey, California 93943 | 1 |
| 16. | Joon Yong Lee
Department of Korean Language
Defense Language Institute
Presidio of Monterey, California 93940 | 1 |
| 17. | Academic Library
Defense Language Institute
Presidio of Monterey, California 93940 | 1 |
| 18. | Professor Young S. Shin, code 69 Sg
Department of Mechanical Engineering
Naval Postgraduate School
Monterey, California 93943 | 1 |

END

FILMED

5-85

DTIC